

# Detecting macroecological patterns in bacterial communities across independent studies of global soils

**Authors:** Kelly S Ramirez<sup>\*,1</sup>, Christopher G. Knight<sup>+2</sup>, Mattias de Hollander<sup>1</sup>, Francis Q. Brearley<sup>3</sup>, Bede Constantinides<sup>4</sup>, Anne Cotton<sup>5</sup>, Si Creer<sup>6</sup>, Thomas W. Crowther<sup>1,7</sup>, John Davison<sup>8</sup>, Manuel Delgado-Baquerizo<sup>9</sup>, Ellen Dorrepaal<sup>10</sup>, David R. Elliott<sup>3,11</sup>, Graeme Fox<sup>3</sup>, Rob Griffiths<sup>12</sup>, Chris Hale<sup>13</sup>, Kyle Hartman<sup>14</sup>, Ashley Houlden<sup>15</sup>, David L. Jones<sup>6</sup>, Eveline J. Krab<sup>10</sup>, Fernando T. Maestre<sup>16</sup>, Krista L. McGuire<sup>17</sup>, Sylvain Monteux<sup>10</sup>, Caroline H. Orr<sup>18</sup>, Wim H van der Putten<sup>1,19</sup>, Ian S. Roberts<sup>15</sup>, David A. Robinson<sup>20</sup>, Jennifer D. Rocca<sup>21</sup>, Jennifer Rowntree<sup>3</sup>, Klaus Schlaeppi<sup>14</sup>, Matthew Shepherd<sup>22</sup>, Brajesh K. Singh<sup>23</sup>, Angela L. Straathof<sup>2</sup>, Jennifer M. Bhatnagar<sup>24</sup>, Cécile Thion<sup>25</sup>, Marcel G.A. van der Heijden<sup>14,26,27</sup>, and Franciska T. de Vries<sup>2</sup>

\* email: k.ramirez@nioo.knaw.nl

+ Joint lead authors

<sup>1</sup> Netherlands Institute of Ecology, Droevendaalsesteeg 10 6708 PB Wageningen NL

<sup>2</sup> Faculty of Science and Engineering, The University of Manchester, Manchester, M13 9PT, United Kingdom.

<sup>3</sup> School of Science and the Environment, Manchester Metropolitan University, Chester Street, Manchester, M1 5GD.

<sup>4</sup> Evolution and Genomic Sciences, School of Biological Sciences, The University of Manchester, Manchester, M13 9PT, United Kingdom

<sup>5</sup> Department of Animal and Plant Sciences, The University of Sheffield, Alfred Denny building, Sheffield, South Yorkshire, S10 2TN, UK

<sup>6</sup> Environment Centre Wales, College of Natural Sciences, Bangor University, Gwynedd, LL57 2UW, United Kingdom.

<sup>7</sup> Institute of Integrative Biology, ETH Zurich, Universitätsstrasse 16, 8006, Zürich, Switzerland.

<sup>8</sup> Department of Botany, Institute of Ecology and Earth Sciences, University of Tartu, Lai 40, Tartu 51005, Estonia

<sup>9</sup> Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder, CO 80309.

<sup>10</sup> Climate Impacts Research Centre, Department of Ecology and Environmental Science, Umeå University, Vetenskapens väg 38, 981 07, Abisko, Sweden

<sup>11</sup> Environmental Sustainability Research Centre, University of Derby, Kedleston Road, Derby, DE22 1GB, UK

<sup>12</sup> Centre for Ecology and Hydrology, Wallingford, United Kingdom

<sup>123</sup> School of Life Sciences, University of Warwick, Coventry, CV4 7AL, United Kingdom

- <sup>14</sup>Division of Agroecology and Environment, Agroscope, Zurich, Reckenholzstrasse 191, Switzerland
- <sup>15</sup>Faculty of Biology, Medicine and Health, The University of Manchester, Manchester, M13 9PT, United Kingdom.
- <sup>16</sup>Departamento de Biología y Geología, Física y Química Inorgánica, Escuela Superior de Ciencias Experimentales y Tecnología, Universidad Rey Juan Carlos, Calle Tulipán s/n, 28933 Móstoles, Spain
- <sup>17</sup>Department of Biology, Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403, USA
- <sup>18</sup>School of Science and Engineering, Teesside University, Middlesbrough, TS1 3BX, United Kingdom.
- <sup>19</sup>Laboratory of Nematology, Wageningen University, Droevendaalsesteeg 1, Wageningen 6708 PB, The Netherlands.
- <sup>20</sup>Centre for Ecology and Hydrology, Bangor, LL57 2UW, United Kingdom
- <sup>21</sup>Department of Biology, Duke University, Durham, NC, 27705, United States
- <sup>22</sup>Natural England, United Kingdom
- <sup>23</sup>Hawkesbury Institute for the Environment, Western Sydney University, Richmond 2753 NSW Australia
- <sup>24</sup>Department of Biology, Boston University, Boston, MA, 02215, United States
- <sup>25</sup>Institute of Biological and Environmental Sciences, University of Aberdeen, Saint-Machar Drive, AB24 3UU, Aberdeen, United Kingdom
- <sup>26</sup>Institute for Evolutionary Biology and Environmental Studies, University of Zürich, Winterthurerstrasse 190, CH-8057, Switzerland.
- <sup>27</sup>Plant-Microbe Interactions, Institute of Environmental Biology, Faculty of Science, Utrecht, The Netherlands

**Keywords:** microbial ecology, soil, diversity, community structure, Illumina sequencing, 16S rRNA gene, biogeography, microbiology, meta-analysis

**The emergence of high-throughput DNA sequencing methods provides unprecedented opportunities to further unravel bacterial biodiversity and its worldwide role from human health to ecosystem functioning. However, in spite of the abundance of sequencing studies, combining data from multiple individual studies to address macroecological questions of bacterial diversity remains methodically challenging and plagued with biases. Here, using a machine learning approach that accounts for differences among studies and complex interactions among taxa, we merge 30 independent bacterial datasets consisting of 1,998 soil samples from across 21 countries. While previous meta-analysis efforts have focused on**

bacterial diversity measures or abundances of major taxa, we show that disparate amplicon sequence data can be combined at the taxonomy-based level to assess bacterial community structure. We find that rarer taxa are more important for structuring soil communities than abundant taxa, and that these rarer taxa are better predictors of community structure than environmental factors, which are often confounded across studies. We conclude that combining data from independent studies can be used to explore bacterial community dynamics, identify potential ‘indicator’ taxa with an important role in structuring communities, and propose hypotheses on the factors that shape bacterial biogeography previously overlooked.

Soil microbial communities are more diverse and contain more individuals than any species groups on the planet<sup>1,2</sup>. Over the last decade, the use of high-throughput sequencing (HTS) methods has substantially advanced our understanding of the worldwide biogeography and ecology of soil bacterial and fungal communities<sup>3–5</sup>. Recent work has further demonstrated that inclusion of microbial composition and functional attributes improves earth system models<sup>6</sup>, which is of paramount importance for predicting effects of global change on ecosystem services such as climate regulation or soil fertility<sup>7</sup>. Yet, opposite to the long-standing view that every organism may occur everywhere<sup>8</sup>, even at small scales bacterial communities turn out to be more patchy than previously expected<sup>9,10</sup>, raising questions regarding dispersal constraints, temporal dynamics, and niche breadth at the global scale<sup>11–13</sup>. Due to these knowledge gaps, combined with practical challenges of exhaustive sample collection and the massive diversity of communities, global assessment of soil microbial diversity remains an ongoing research challenge<sup>14</sup>.

82

83 For plants and animals, the integration of data from independent studies has been a valuable  
84 option for generating an understanding of global biogeography patterns, answering ecological  
85 questions (e.g. biodiversity-functioning relationships), and identifying threats to biodiversity  
86 from global changes<sup>15-17</sup>. Similarly, our understanding of soil microbial diversity would greatly  
87 improve from such worldwide assessments. However, the integration of microbial community  
88 HTS data from different studies is not so unlike the merging of museum species records where  
89 information and data is constrained by variations in nomenclature over space and time, among  
90 many other challenges<sup>18,19</sup>. Like plant and animal records, molecular microbial community  
91 records and information can be incomplete, processing and naming varies greatly between  
92 studies and over time<sup>20</sup>, data storage is inconsistent, and there are few curated databases with  
93 high quality data (especially for short read sequences)<sup>21,22</sup>. Further, most microbial community  
94 data and metadata are still available only in independently published studies that have been  
95 carried out according to their own standards and procedures, and the extent of these confounding  
96 factors has never been quantified across studies.

97

98 Regardless of the challenges, as indicated by the many open access data initiatives<sup>23-25</sup>, merging  
99 microbial sequence data is a potential option to address global scale questions, whether relating  
100 to the human microbiome<sup>26</sup>, marine systems<sup>27</sup>, or predicting the response of soil organisms to  
101 global environmental change<sup>28</sup>. For soil systems, the need to merge sequence data is supported  
102 by the emerging role of bacterial phyla and classes as indicators of particular soil conditions such  
103 as soil pH and nutrient concentrations<sup>29,30</sup>. Until now, attempts to meta-analyze sequence data  
104 have been limited to assessing diversity measures or abundances of major taxa, because the

merging of community data is constrained by methodological differences between sequencing studies<sup>10,24,31,32</sup>. However, a recent systematic review found that measures of microbial community structure were more often linked to microbial process rates than diversity or presence/absence data<sup>33</sup>, and abundance ratios among phyla may be less important than previously believed<sup>34</sup>. Together indicating that information on variation in microbial community structure is potentially more ecologically relevant than measures of diversity and abundances of major taxa.

Here, we show that, despite the outlined challenges, published microbial community data from independent studies can be analyzed together to address questions about the global structuring of communities. Using a machine learning approach, we take methodological and technical biases into account, factor in interactions among taxa, and produce an improved assessment of the abiotic and biotic drivers of soil community structure. The objectives of this study were two-fold: (1) to identify the biases and incompatibilities of microbial community HTS studies (and confounding factors) so as to strengthen our ability to integrate data from disparate studies, and (2) to reveal worldwide soil microbial community patterns by merging independent taxonomy-based datasets.

## **Results and Discussion**

### **Taxonomy-based merging of disparate amplicon sequence data**

We identified 30 individual HTS bacterial studies from 21 countries for our analysis (Figure 1A and Supplementary Table 1). While we aimed to merge HTS data of both soil bacterial and fungal datasets, our approach was only successful for bacterial data (Figure 1B and 1C), and highlights the well-known dilemma of fungal databases, where extremely high diversity

combined with high endemism and mismatched taxonomy across continents make merging data by taxonomy difficult and unusable for downstream analyses<sup>4,35</sup>. For the bacterial studies, we were able to successfully merge 30 individual OTU tables; using a taxonomy-based approach, datasets were merged using the taxonomic affiliations of individual OTUs. Once filtered, and singletons removed, the final ‘taxonomy-based’ community contained 1,998 individual soil samples, and 8,287 taxa. Here ‘taxon’ is defined as a unique name in the classification; where a name could be a specific phylum, genus, or other taxonomic level. For example, ‘Acidovorax’ (genus) and Proteobacteria (the phylum containing Acidovorax) were both considered as taxa). To account for variation in sequencing depth between different studies, OTU relative abundances were used per sample, rather than absolute read abundance. To test known biogeographical patterns, metadata (information on geographical location, soil pH and soil core measurements) were compiled for all studies. Technical and methodical information was also collected; all of these 30 studies had conducted amplicon sequencing on hypervariable regions of the 16S rRNA gene in soil samples using either Illumina or (Roche) 454 pyrosequencing (with any primer pair) (Supplementary Table 1). For a validation step we retrieved all usable raw sequence data available, resulting in 417 samples from locations across the globe (approximately 1/5 of all our samples) (Figure 1A). Data not included in this sequence-matched analysis either had an incompatible raw sequence format or simply no longer existed. Available raw sequence data were combined into a single ‘sequence-matched’ community comprising 44,106 OTUs (Supplementary Figure 1).

## **Machine learning assessment of bacterial community structure**

Ordination of the taxonomy-based community reveals large amounts of structure both within and

between studies (structure that is removed by permuting taxa among samples (Supplementary Figure 2), without greatly affecting diversity (Supplementary Table 3)), and the observation of the well-established negative relationship between relative abundance of Acidobacteria and soil pH (Figure 1D)<sup>36</sup> confirms our merging method. This visualization also suggests that some of the community variation (e.g. the near absence of Acidobacteria in some studies, even at low pH) is due to technical factors such as the particular primer sets chosen, region sequenced, and sequencing platform (Supplementary Methods and Supplementary Table 2). However, we expect that some taxa are not correlated with technical factors, and are non-randomly distributed with respect to biotic and abiotic factors. Therefore, using a machine learning approach capable of accounting for complex interactions among taxa (Random Forests<sup>TM</sup>, see methods), we determined the extent to which individual taxa could influence the community structure of merged independent studies. Here community structure is defined by the presence and relative abundances of individual taxa, along with co-occurrence relationships between those taxa. This was done in two ways: first, we constructed a model that classified the study from which a sample came based on the proportions of the 8,287 taxa it contained (1.5% [ $\pm$  0.02% CI] classification error, by internal cross-validation). Second, we determined the contribution of each taxon to bacterial community structure by quantifying its importance in a model that separated the observed data from synthetic data randomly drawn from the observed distributions of relative abundances for each taxon (*see Methods*).

Merging of disparate microbial sequence data is known to be plagued with potential biases including: lack of standardization of sample collection, methodological issues regarding DNA extraction and primer choice, incomplete metadata, the technical biases of different sequencing

platforms, sequencing depth, PCR Bias, different clustering methods, and the use of different taxonomic classification pipelines<sup>37–39</sup>. We therefore took the step to quantify the importance of both technical and environmental factors alongside taxa in the Random Forests models (Figure 2). Of note, ‘owner’, which encompasses the technical biases and uniqueness of a given dataset, is very effective for differentiating between studies (i.e. the owner is far to the right in Figure 2) yet is entirely uninformative about community structure (i.e. owner is at the far bottom in Figure 2). In fact, *all* technical factors included are better than 98.5% of all taxa to differentiate between studies, indicating that the observed differences among studies in taxon relative abundances are strongly confounded with technical factors. Independent of taxonomy, certain environmental factors, such as country of origin, latitude and longitude, and soil pH, were highly important in differentiating studies but not in determining community structure. By contrast, minimum soil sampling depth was not very important in separating studies, and was more associated with community structure. It is well known that bacterial diversity decreases with soil depth<sup>40</sup> and our results show that in a global assessment, soil depth remains a strong predictor of bacterial community composition. Perhaps most useful for future research, this result highlights that not all environmental factors are equally confounded by technical factors, and shows that by combining data from across many independent studies we may identify previously overlooked taxa and factors relevant for structuring communities.

### **Importance for structuring soil bacterial communities**

Although all studies were confounded by technical and environmental covariates, there remained many taxa that were non-randomly distributed and were not confounded with technical differences among studies (upper left in Figure 2). When assessing the role of these different taxa



in structuring the community, we found a trade-off between taxon abundance and importance in community structure, such that low abundance taxa are disproportionately important in the non-random structure of communities, where the most important taxa are rarer than expected compared to the randomly permuted data (Figure 3). Thus, the importance of taxa for determining community structure is negatively correlated with the average abundance of those taxa, whereas taxon abundance is positively correlated with importance for separating studies ( $\rho = -0.79$  and  $\rho = +0.51$  respectively, rank correlation, cf. null expectations of  $\rho = -0.62$  and  $-0.12$  respectively in permuted data). The taxa most closely associated with differences between studies tend to be those present at or greater than 0.1% relative abundance, but those most important in determining community structure tend to be present at 0.0001% abundance or less (with a null expectation of around 0.01-0.001% in each case, Figure 3). This result is only found by considering the full set of studies and is neither apparent within single studies (Supplementary Fig. 4A-B) nor a subset of studies (whether matched by name or sequence Supplementary Fig. 5). It corresponds to the long tail in frequency-abundance distributions of soil microbial communities<sup>41</sup>, where many taxa in the soil are known to occur at low abundance. Thus if rarer taxa tend to be more important for distinguishing between communities, it is within this long tail that we might identify taxa that could indicate ecological or functional differences among soil communities<sup>42,43</sup>.

To be ecological indicators<sup>44,45</sup>, taxa need to vary in abundance in response to environmental factors and have high occurrence across studies, as is the case for the phylum Acidobacteria<sup>36</sup>. Acidobacteria, however, are typically abundant and our analysis suggests that the most abundant taxa are *not* the most important in determining community structure. While dominant taxa like

Acidobacteria do change with environmental factors such as pH (Figure 1D), those changes are of lesser importance for the ‘non-randomness’ of community structure, and more confounded with technical effects, than changes in less dominant, pH responsive taxa (Supplementary Figure 3A). Therefore, we assessed which taxonomic ranks are more or less distinguished from the randomly permuted data. Although differences among domains and phyla are strongly associated with differences among studies (Figure 4B) only taxa at a rank lower than phyla are consistently better than random at identifying community structure (Figure 4A).

A very similar pattern was found for the sequence-matched community, emphasizing the importance of taxa at the level of Class and below (Supplementary Figure 7A and 7B). However, this was not apparent in individual studies (Supplementary Figure 4C-D), where phyla were relatively important. A subset of the taxonomy-matched studies showed a pattern intermediate between the single studies and the full dataset (phyla with some importance, but less than Class, Order or Family, Supplementary Figure 7C). This, along with abundance analyses (Figure 3 and Supplementary Figure 5), suggests that our name matching approach is consistent with, but less powerful than a full sequence-matched analysis. At the same time, the taxonomy-matching is worthwhile because, as with the findings on abundance (Figure 3), macroecological patterns (the importance of taxa below phyla and of relatively low abundance in community structure) are evident when we consider thousands of samples from tens of studies, that are not apparent from hundreds of samples from one or a handful of studies.

To be a good ecological indicator a taxon should occur in most studies; we therefore looked explicitly at the relationship between a taxon’s importance in community structure and its

occurrence across studies. Low abundance taxa and taxa of lower taxonomic rank are consistently important in determining community structure, but tend to be detected in fewer studies ( $\rho = 0.59$  and  $0.31$  respectively Supplementary Figure 3B and 3C). We discovered a relationship between taxon occurrence across studies and importance for structuring communities for all taxa (Figure 5, Supplementary Table 4). Comparison with the null expectation reveals a range of taxa, occurring in multiple samples from most studies, which are much more important in determining community structure than expected by chance. A similar pattern is apparent in the sequence-matched dataset (Supplementary Figure 8A) and the same subset of studies when taxonomy-matched (Supplementary Figure 8B). Altogether, the analysis clearly illustrates the significance of taxonomic rank, for example *class* Gemmatimonadetes is relatively unimportant for community structure but *genus* Gemmatimonadetes is relatively important. The result also shows rarer taxa being more important in structuring communities and suggests rarer bacterial taxa play overlooked ecologically important roles for bacterial community dynamics<sup>43</sup>. This result is robust to artifacts caused by the rarest taxa (e.g. differences between 0 and 1 reads in a sample could be significant for a model, without being biologically significant) – a very similar pattern is seen when only taxa present at above 0.003% in any given sample were included in this analysis (typically removing the rarest 10% of taxa from any given sample, Supplementary Figure 9). Conversely, many taxa of high taxonomic rank with high occurrence across samples, such as the phyla Actinobacteria, Acidobacteria, Proteobacteria, and Bacteroidetes, were much less important for community structure than the null expectation. These taxa have been reported elsewhere as ‘core’ members of the soil community<sup>36,46</sup>, and even been included in source-tracking of microbial communities due to their ubiquitous presence in soil<sup>47</sup>. Yet, it is the consistent presence of the core taxa across samples

and studies that makes them inadequate for assessing community structure.

## **Conclusions**

Our results demonstrate the power of combining global bacterial HTS data from multiple independent sources for the detection of biogeographical patterns and for identifying community patterns that can be used to generate hypotheses on the roles of certain taxa. Though our assessment was on soil communities, our methods can be applied to broadly to other microbial datasets and disciplines. Taxonomy-based merging gives results that are consistent with raw sequence data, and expands opportunities for extracting information about microbial communities from the wealth of existing and future studies. Moreover, we find that rarer bacterial taxa are more important in differentiating communities than previously assumed, and hold potential as overlooked soil indicators or keystone species. Still, there are considerable challenges associated with merging large sequence datasets beyond the well-known biases that accompany any molecular HTS study. Perhaps the most concerning was that so few raw sequence datasets for publically deposited analyses could be retrieved. This highlights the need for wider community adoption of open and accessible short read sequence databases<sup>48</sup>, open reference clustering<sup>49</sup>, standardized databases<sup>50</sup> and—as always—that metadata should be consistent and accessible. Regardless of these challenges, as HTS methods rapidly advance we must find ways to simultaneously curate and carry our research knowledge forward. Only then, in combination with the many recently designed and classical approaches, can we uncover the full breadth of soil diversity and the roles soil microbes play for ecosystem processes.

## Methods:

### *Description of datasets:*

Metadata from the 30 studies and 1998 samples were collected and compiled into a summary data file. To do so, we standardized the metadata of each study using the dplyr package<sup>51</sup> of the R statistical platform<sup>52</sup>. Samples were collected from 21 countries representing all continents except Antarctica. In addition to location and pH data (median = 6.1, quartile range=5.3-7.0), which were available from all studies, information on altitude (10 m, 10-860 m), soil moisture (19.5%, 14.1-27.4%), and total soil nitrogen (0.36 mg kg<sup>-1</sup>, 0.23-0.51 mg kg<sup>-1</sup>), carbon (4.7%, 1.9-7.5%) and phosphorus (20.7 mg kg<sup>-1</sup>, 7.0-223.0 mg kg<sup>-1</sup>) was noted where available. Depth of sample collection was also noted and ranged from surface collections to a maximum depth of 70 cm, with 83% of samples originating from 0-10 cm below the soil surface. Samples represented anthropogenically managed (59%) and natural (40%; remaining samples undefined) systems, and were taken from arable, grassland, peatland, forest, scrub (including tundra) and urban habitats. The majority of samples (71%) were described as non-experimental, meaning no treatments were applied, with the remainder described as experimental. Sequencing data were either produced using Roche 454 technology (22%) or one of the Illumina platforms (78%). Primer pairs were defined for 92% of the samples and nine different pairs were identified from the study meta data (27F:338R; 341F:518R; 341F:806R; 341F:907R; 357F:926R; 515F:806R; 577F:926R; 799F:1193R and 341F:805R) with the majority of samples (66%) using 515F and 806R to produce amplicons. Post sequencing processing varied, but 81% of samples were run through the QIIME workflow at some point. An OTU table for 1 study comprising 43 samples was programmatically retrieved from the MG-RAST public metagenome repository<sup>53</sup>. Taxonomy for the different studies was mainly assigned using the Greengenes database (84 %),

but RDP (6 %;<sup>37</sup> and the Silva database (9 %)<sup>54</sup> were also used.

### *Primer Biases*

It has long been well understood that different primers vary in their biases for amplifying members of the bacterial community<sup>55,56</sup>. To demonstrate this bias, the likelihood of significant differences in primer biases for the ten pairs of primers used in the studies analysed were determined by *in silico* analysis. Sequences of primer pairs were compared to all 16S rRNA gene sequences in the SILVA non-redundant reference database (SSURef NR) release 128<sup>54</sup> using TestPrime v1.0 (as described in<sup>57</sup>). The percentages of sequences of each bacterial phyla that matched both primers (with a one base pair mismatch allowance at least 1bp from the 3' end of the primers) were calculated to compare predicted differences in primer coverage of different bacterial taxa.

### *Merging OTU tables:*

For the OTU tables from the 30 individual studies to be merged, extensive data cleaning was carried out on the OTU and taxonomy files to maximize the possibility of matching taxa across datasets. This comprised several steps: (1) Most datasets contained a seven-level taxonomy, recorded in a variety of ways, which was converted to a standardized format. (2) Individual taxon names were cleaned, to give a single name at each taxonomic level (e.g. removing special characters and extra annotations, such as 'candidate division' or details of containing taxa). (3) For the many cases where a taxon was not assigned at a particular taxonomic level, a unified 'unassigned' label was created. Repeating analyses with all these taxa removed made no qualitative difference to the results (Supplementary Figure 10). Merging at the taxonomy-based

level has the added benefit of lessening the impacts of hypervariable regions. For example, the identification of an organism at a specific level in one sample also contributes to the identification of the containing genus for that sample, allowing direct comparison with a sample where, because a different region was sequenced, that same organism is only resolved to the genus level. Next, relative abundance data were, where necessary, re-scaled to sum to 1 for a sample, using original OTU count files where possible. These values were then manipulated to give data tables usable for modeling using custom R scripts. For some analyses (Figures 3-5), a dataset without community structure was created by randomly permuting the relative abundance of each taxon across all samples. Unless otherwise stated, the analyses performed on the permuted dataset was identical to that performed on the observed data.

#### *Merging raw sequence data and other validation datasets:*

While no dataset can currently provide a “ground truth” against which to judge our approach, we can at least validate it. The primary validation of our taxonomy-matching approach was to merge raw sequence data (‘sequence-matched’) from 419 samples of the total 1998 used. Per sample fastq files were obtained for each individual dataset. Read files were quality filtered with sickle<sup>58</sup> for single end reads trimming bases below phred score 36 and shorter than 100bp. These stringent filtering criteria were applied to keep only high quality reads and to make sure it is possible to map reads to full length 16S rRNA gene sequences. Full length 16S rRNA gene sequences from the Silva 119 release<sup>54</sup> were obtained in Qiime compatible format from the [Silva Download Archive](#). For each dataset, all reads were mapped to the full length 16S rRNA gene sequences using the usearch global algorithm implemented in VSEARCH version 1.9.6<sup>59</sup>. The alignment results in usearch table format (uc) were directly converted to BIOM format using

biom version 2.1.5<sup>60</sup>. Consensus/majority taxonomy was added as metadata to the biom file. Finally, all BIOM files of each dataset were merged using Qiime version 1.9.1<sup>61</sup>. All steps were implemented in a workflow made with Snakemake version 3.5.4<sup>62</sup> available: ([De Hollander 2016](#)) (Supplementary Figure 1).

To use this sequence-matched dataset to validate our taxonomy-matching approach across studies using different taxonomy databases (Supplementary Figures 5, 7 & 8) we created an equivalent taxonomy-matched dataset from the same 5 studies. As with the full dataset, only taxa occurring in at least two studies were included in either this or the sequence-matched dataset. To test what is gained or lost by considering different numbers of studies simultaneously, we considered, not only the full dataset (30 studies) and the subset of 5 studies used in the sequence-matched dataset, but two of the largest individual studies: from Central Park, NYC encompassing 594 samples (study #24) and a global dataset encompassing 103 samples (study #30). In each case a simple subset of the full dataset was analyzed (Supplementary Figure 4). To address PCR biases (Supplementary Table 2) and biases associated with rare taxa, we created a filtered subset of the data where only taxa present at above 0.003% in any given sample were considered, meaning that all taxa deemed present are represented by multiple sequence reads (Supplementary Figure 9). To address the issue of differential 16S copy numbers skewing abundance estimates, we created a binary dataset of the presence/absence of all taxa. The results for a model separating studies using this dataset were very similar to the main dataset using relative abundance, however, there was insufficient power to identify taxa important for community structure. Nonetheless, this analysis did agree with the main analysis that phyla were the most stable taxonomic level, with lower importance than on the permuted data



(Supplementary Figure 6). Finally, to test the effect of ‘unknown’ or unclassified bacterial taxa we created a reduced dataset where all taxa classified as ‘unassigned’ at any level were removed (Supplementary Figure 10).

#### *Random forest models.*

To test for the importance of different taxa in the structuring of the data we used Random Forest models<sup>63–65</sup> with the relative abundances of the taxa as explanatory variables. Random Forest models have two principal advantages in this context: 1) they can deal easily with thousands of explanatory variables and quantify their relative importance, and 2) they can run equivalently in both supervised and un-supervised modes. In the latter, the importance of a variable describes how effective it is at separating the observed data from randomized synthetic data<sup>65</sup>. In both cases, a proximity matrix may be generated, which can be used for ordination (Supplementary Figure 2). The importance of individual taxa in a Random Forest relate to traditional ecological measures. For instance, the importance in a supervised model, such as that used separating studies (x-axis in Figure 2) is closely correlated with the sensitivity component of the indicator value of each taxon ( $\rho = 0.89$ , Supplementary Figure 3D)<sup>45</sup>. There are two key parameters that may be adjusted in a Random Forest model, *mtry*, the number of variables randomly sampled as candidates for a split in the constituent trees and *ntree*, the number of trees in the forest. *mtry* was set at its default value (square root of the number of variables) *ntree* was set to 100,000 for each forest. Such a large number of trees was found to be necessary to achieve stable importance across taxa and was achieved by combining several forests run in parallel without normalizing votes. Other parameters were left at default values, in particular, trees were grown to completion (i.e. a minimum node size of 1). The un-scaled permutation importance of variables is used

throughout: Each variable importance is the difference between the classification error rate of a tree on data not used to construct it (the ‘out of bag’ data) and the same error following random permutation of the variable in question, averaged over all trees.

We used permuted data (see above) to create null distributions for taxon importances. For unsupervised Random Forests analyses, such as the community structure model, this amounts to calculating how important a taxon with a particular abundance distribution is for separating two randomized distributions. This can then be compared to its importance for separating the observed from a randomized distribution. This clarifies the fact that, even in null data without community structure (Supplementary Figure 2), variable importance correlates with ecologically important factors, such as abundance. This makes intuitive sense in as much as, even with randomized samples, is easier to separate them on the basis of taxa that occur in only some of them than on the basis of ubiquitous taxa. This, for instance, results in the negative slope of the orange (permuted, null, data) line in Figure 5. All analyses were completed with RandomForest package for R version 4.6.

## References

1. Prosser, J. I. Dispersing misconceptions and identifying opportunities for the use of ‘omics’ in soil microbial ecology. *Nat. Rev. Microbiol.* **13**, 439–46 (2015).
2. Bardgett, R. D. & van der Putten, W. H. Belowground biodiversity and ecosystem functioning. *Nature* **515**, 505–511 (2014).
3. Caporaso, J. G. *et al.* Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J.* **6**, 1621–4 (2012).
4. Tedersoo, L. *et al.* Global diversity and geography of soil fungi. *Science (80-. ).* **346**, (2014).
5. Davison, J. *et al.* Global assessment of arbuscular mycorrhizal fungus diversity reveals very low endemism. *Science (80-. ).* **349**, (2015).
6. Wieder, W. R., Bonan, G. B. & Allison, S. D. Global soil carbon projections are improved by modelling microbial processes. *Nat. Clim. Chang.* **3**, 909–912 (2013).

7. Karhu, K. *et al.* Temperature sensitivity of soil respiration rates enhanced by microbial community response. *Nature* **513**, 81–84 (2014).
8. Barberán, A., Casamayor, E. O. & Fierer, N. The microbial contribution to macroecology. *Front. Microbiol.* **5**, 203 (2014).
9. Ramirez, K. S. *et al.* Biogeographic patterns in below-ground diversity in New York City's Central Park are similar to those observed globally. *Proc. Biol. Sci.* **281**, 20141988- (2014).
10. O'Brien, S. L. *et al.* Spatial scale drives patterns in soil bacterial diversity. *Environ. Microbiol.* **18**, 2039–2051 (2016).
11. Evans, S., Martiny, J. B. H. & Allison, S. D. Effects of dispersal and selection on stochastic assembly in microbial communities. *ISME J.* **11**, 176–185 (2017).
12. Talbot, J. M. *et al.* Endemism and functional convergence across the North American soil mycobiome. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 6341–6 (2014).
13. Barber, A. *et al.* Why are some microbes more ubiquitous than others? Predicting the habitat breadth of soil bacteria. *Ecol. Lett.* **17**, 794–802 (2014).
14. Ranjard, L. *et al.* Turnover of soil bacterial diversity driven by wide-scale environmental heterogeneity. *Nat. Commun.* **4**, 1434 (2013).
15. Jetz, W., McPherson, J. M. & Guralnick, R. P. Integrating biodiversity distribution knowledge: toward a global map of life. *Trends Ecol. Evol.* **27**, 151–9 (2012).
16. Ricketts, T. H. *et al.* Disaggregating the evidence linking biodiversity and ecosystem services. *Nat. Commun.* **7**, 13106 (2016).
17. Dirzo, R. *et al.* Defaunation in the Anthropocene. *Science (80-. )*. **345**, 401–406 (2014).
18. Patterson, D. J., Cooper, J., Kirk, P. M., Pyle, R. L. & Remsen, D. P. Names are key to the big new biology. *Trends Ecol. Evol.* **25**, 686–691 (2010).
19. Santos, A. M. & Branco, M. The quality of name-based species records in databases. *Trends Ecol. Evol.* **27**, 6-7-8 (2012).
20. Beiko, R. G. Microbial Malaise: How Can We Classify the Microbiome? *Trends Microbiol.* **23**, 671–679 (2015).
21. Tedersoo, L. *et al.* Standardizing metadata and taxonomic identification in metabarcoding studies. *Gigascience* **4**, 34 (2015).
22. Ramirez, K. S. *et al.* Toward a global platform for linking soil biodiversity data. *Front. Ecol. Evol.* **3**, (2015).
23. Turner, W. *et al.* Free and open-access satellite data are key to biodiversity conservation. *Biol. Conserv.* **182**, 173–176 (2015).
24. Gilbert, J. A., Jansson, J. K. & Knight, R. The Earth Microbiome project: successes and aspirations. *BMC Biol.* **12**, 69 (2014).
25. Joppa, L. N. *et al.* Filling in biodiversity threat gaps. *Science (80-. )*. **352**, (2016).
26. Sinha, R., Abnet, C. C., White, O., Knight, R. & Huttenhower, C. The microbiome quality control project: baseline study design and future directions. *Genome Biol.* **16**, 276 (2015).
27. Sogin, M. L. *et al.* Microbial diversity in the deep sea and the underexplored 'rare biosphere'. *Proc. Natl. Acad. Sci. U. S. A.* **103**, 12115–20 (2006).
28. García-Palacios, P. *et al.* Are there links between responses of soil microbes and ecosystem functioning to elevated CO<sub>2</sub>, N deposition and warming? A global perspective. *Glob. Chang. Biol.* **21**, 1590–1600 (2015).
29. Hermans, S. M. *et al.* Bacteria as Emerging Indicators of Soil Condition. *Appl. Environ. Microbiol.* **83**, AEM.02826-16 (2017).

- 478 30. Philippot, L. *et al.* The ecological coherence of high bacterial taxonomic ranks. *Nat. Rev.*  
479 *Microbiol.* **8**, 523–529 (2010).
- 480 31. Shade, A., Caporaso, J. G., Handelsman, J., Knight, R. & Fierer, N. A meta-analysis of  
481 changes in bacterial and archaeal communities with time. *ISME J.* **7**, 1493–506 (2013).
- 482 32. Hendershot, J. N., Read, Q. D., Henning, J. A., Sanders, N. J. & Classen, A. T.  
483 Consistently inconsistent drivers of microbial diversity and abundance at macroecological  
484 scales. *Ecology* **98**, 1757–1763 (2017).
- 485 33. Bier, R. L. *et al.* Linking microbial community structure and microbial processes: an  
486 empirical and conceptual overview. *FEMS Microbiol. Ecol.* **91**, (2015).
- 487 34. Walters, W. A., Xu, Z. & Knight, R. Meta-analyses of human gut microbes associated  
488 with obesity and IBD. *FEBS Lett.* **588**, 4223–4233 (2014).
- 489 35. Bik, H. M. *et al.* Sequencing our way towards understanding global eukaryotic  
490 biodiversity. *Trends Ecol. Evol.* **27**, 233–243 (2012).
- 491 36. Lauber, C. L., Hamady, M., Knight, R. & Fierer, N. Pyrosequencing-Based Assessment of  
492 Soil pH as a Predictor of Soil Bacterial Community Structure at the Continental Scale.  
493 *Appl. Environ. Microbiol.* **75**, 5111–5120 (2009).
- 494 37. McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological  
495 and evolutionary analyses of bacteria and archaea. *ISME J.* **6**, 610–8 (2012).
- 496 38. Lozupone, C. & Stombaugh, J. Meta-analyses of studies of the human microbiota.  
497 *Genome ...* (2013).
- 498 39. Pawluczyk, M. *et al.* Quantitative evaluation of bias in PCR amplification and next-  
499 generation sequencing derived from metabarcoding samples. *Anal. Bioanal. Chem.* **407**,  
500 1841–1848 (2015).
- 501 40. Lu, X., Seuradze, B. J. & Neufeld, J. D. Biogeography of soil Thaumarchaeota in relation  
502 to soil depth and land usage. *FEMS Microbiol. Ecol.* **93**, (2017).
- 503 41. Jung, S. P. & Kang, H. Assessment of microbial diversity bias associated with soil  
504 heterogeneity and sequencing resolution in pyrosequencing analyses. *J. Microbiol.* **52**,  
505 574–580 (2014).
- 506 42. Langille, M., Zaneveld, J. & Caporaso, J. Predictive functional profiling of microbial  
507 communities using 16S rRNA marker gene sequences. *Nature* (2013).
- 508 43. Jousset, A. *et al.* Where less may be more: how the rare biosphere pulls ecosystems  
509 strings. *ISME J.* **11**, 853–862 (2017).
- 510 44. Hermans, S. M. *et al.* Bacteria as emerging indicators of soil condition. *Appl. Environ.*  
511 *Microbiol.* AEM.02826-16 (2016). doi:10.1128/AEM.02826-16
- 512 45. Cáceres, M. De & Legendre, P. Associations between species and groups of sites: indices  
513 and statistical inference. *Ecology* **90**, 3566–3574 (2009).
- 514 46. Maestre, F. T. *et al.* Increasing aridity reduces soil microbial diversity and abundance in  
515 global drylands. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 15684–9 (2015).
- 516 47. Knights, D. *et al.* Bayesian community-wide culture-independent microbial source  
517 tracking. *Nat. Methods* **8**, 761–763 (2011).
- 518 48. Muir, P. *et al.* The real cost of sequencing: scaling computation to keep pace with data  
519 generation. *Genome Biol.* **17**, 53 (2016).
- 520 49. Rideout, J. R. *et al.* Subsampled open-reference clustering creates consistent,  
521 comprehensive OTU definitions and scales to billions of sequences. *PeerJ* **2**, e545 (2014).
- 522 50. Yilmaz, P. *et al.* The genomic standards consortium: bringing standards to life for  
523 microbial ecology. *ISME J.* **5**, 1565–7 (2011).

51. Wickham, H. & Francois, R. dplyr: A Grammar of Data Manipulation. R package version 0.5.0. *R package version 0.5.0*. (2016). at <<https://cran.r-project.org/package=dplyr>>
52. Computing., R. A. language and environment for statistical. R Core Team. (2016).
53. Wilke, A. *et al.* The MG-RAST metagenomics database and portal in 2015. *Nucleic Acids Res.* **44**, D590–D594 (2016).
54. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.* **41**, D590–6 (2013).
55. Suzuki, M. T. & Giovannoni, S. J. Bias caused by template annealing in the amplification of mixtures of 16S rRNA genes by PCR. *Appl. Environ. Microbiol.* **62**, 625–30 (1996).
56. Sipos, R. *et al.* Effect of primer mismatch, annealing temperature and PCR cycle number on 16S rRNA gene-targeting bacterial community analysis. *FEMS Microbiol. Ecol.* **60**, 341–350 (2007).
57. Klindworth, A. *et al.* Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res.* **41**, e1 (2013).
58. Joshi & Fass, J. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files. (2011).
59. Rognes, T. *et al.* vsearch: VSEARCH 1.9.6. (2016). doi:10.5281/ZENODO.44512
60. McDonald, D. *et al.* The Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love the ome-ome. *Gigascience* **1**, 7 (2012).
61. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing data. *Nat Meth* **7**, 335–336 (2010).
62. Koster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics* **28**, 2520–2522 (2012).
63. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
64. Breiman, L. & Cutler, A. Random Forests Manual v4.0. *Technical report, UC Berkeley* (2003). at <<https://www.scribd.com/document/208387804/Using-Random-Forests-v4-0>>
65. Shi, T. & Horvath, S. Unsupervised Learning With Random Forest Predictors. *J. Comput. Graph. Stat.* **15**, 118–138 (2006).

**Data availability:** The authors declare that the data supporting the findings of this study are available within the paper and its supplementary information files.

**Correspondence and requests for materials** should be addressed to K.S.R

**Acknowledgements:** We thank all the people who contributed data and input to this study. This study was conducted at a workshop (5/2015, Manchester, UK) funded by the British Ecological Society's special interest group Plants-Soils-Ecosystems and organized by FTDV and KSR. This study and participants were funded in part by ERC Adv grant 26055290 (KSR, WHvdP); BBSRC David Phillips Fellowship (BB/L02456X/1) (FTDV); ERC Grant Agreements 242658 [BIOCOM] and 647038 [BIODESERT] (FTM); the European Regional Development Fund (*Centre of Excellence EcolChange*) (JD); Yorkshire Agricultural Society, Nafferton Ecological Farming Group, and the Northumbria University Research Development Fund (CHO); BBSRC Training Grant (BB/K501943/1) (CH); Wallenberg Academy Fellowship (KAW 2012.0152), Formas (214-2011-788) and Vetenskapsrådet (612-2011-5444) (ED); the Glastir Monitoring & Evaluation Programme (Contract reference: C147/2010/11) and the full support of the GMEP team on the Glastir project (DLJ, SC, DAR). Data taken from work carried out in collaboration with CJS, CL, JMC and SPC. Computing was facilitated by the University of Manchester Condor pool and the CLIMB infrastructure ([www.climb.ac.uk](http://www.climb.ac.uk)).

**Author Contributions:** F.T.dV. and K.S.R. conceived the idea of this study. The datasets were compiled by C.G.K., R.G., J.D., A.H., B.C., G.F., A.L.S. & J.K.R.. Metadata was compiled by J.D and J.K.R.. Raw sequence analysis was conducted by M.dH.. Primer bias analysis was conducted by A.C.. Random forest analyses and figures were conducted by C.G.K.. The manuscript was written by K.S.R., C.G.K., and F.T.dF. with contributions from all co-authors.

**Figures:**

**Figure 1. Merging of data from 32 independent studies demonstrates wide geographic breadth, community variation, and confirms the well-known importance of soil pH. A.** Map of locations from which samples were collected, with zoom panels on the United States (left) and western Europe (right). Points in blue were used in both the taxonomy-based and raw-unified analyses and red points were only used in taxonomy-based analyses. **B.** Average proportion of total prokaryotic abundance and **C.** eukaryotic abundance, represented by taxa shared among different numbers of datasets at different taxonomic levels. Level 1 indicates the complete data, levels 2-4 are subsets of the data containing only taxa present in a minimum of 2-4 separate datasets. **D.** Correlation plot of Acidobacteria relative abundance to soil pH where each color represents a different study ( $r = -0.42$   $p = 8.6 \times 10^{-87}$ ).

**Figure 2: Regardless of technical differences between studies, many bacterial taxa are still informative about bacterial community structure.** Machine learning models classify the study from which samples came (x-axis) based on the relative abundance of taxa within samples and distinguish the observed distribution of taxa among samples from random (y-axis). Plotted alongside bacterial taxa (black) are technical factors (red) and ecological factors (purple), including soil pH, minimum and maximum soil depth, longitude, latitude and degrees from the equator. All values are variable importance from Random Forest models (see *Methods*) – points further to the right on the x-axis have more importance in separating studies, while points higher up on the y-axis, have more importance for community structure. Note the non-linear axes.

**Figure 3: Rarer taxa are more important for structuring communities than abundant taxa.**

Here we show the thousand most important bacterial taxa in community structure (A) and in separating studies (B) with respect to their average relative abundance across samples. Plotted are the ‘observed’ points (green) and ‘permuted’ points (orange) which are a null distribution from performing the same analysis on a permuted dataset (see *Methods*). The y-axis reports the rank variable importance in the Random Forests model of community structure (see *Methods*), i.e. the taxon with the greatest importance in this model is ranked 1, the second greatest 2, etc.

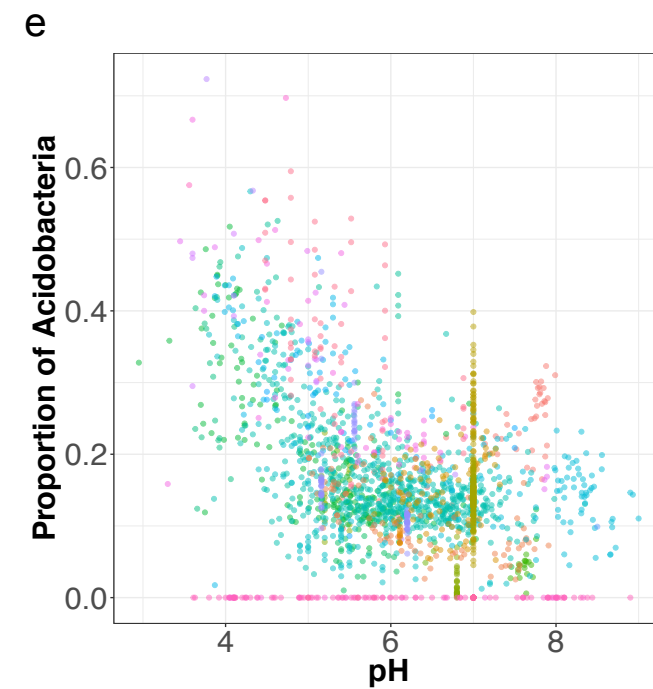
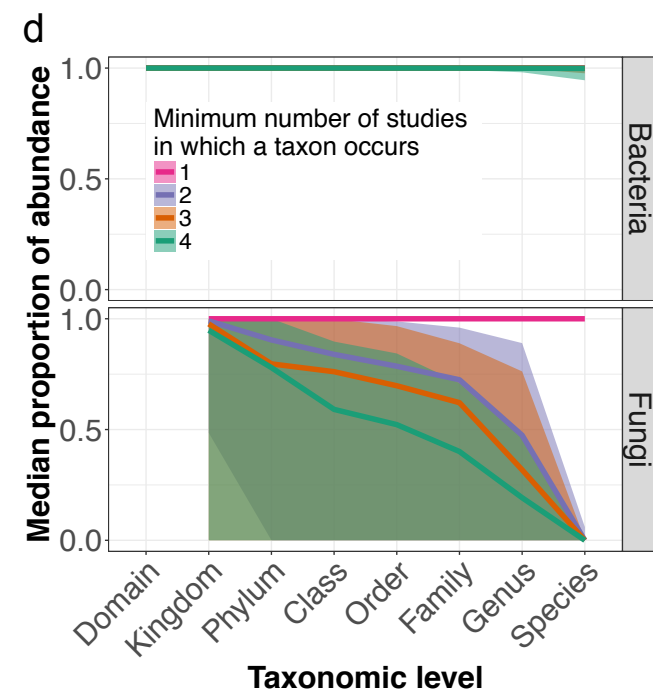
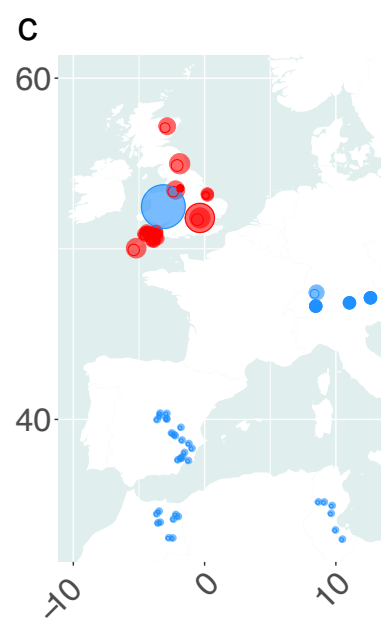
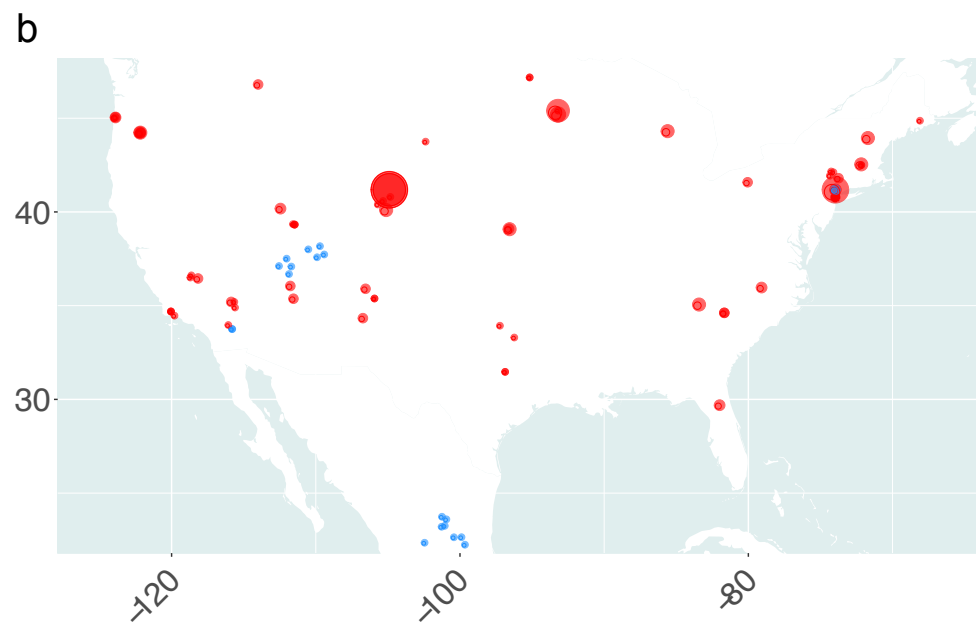
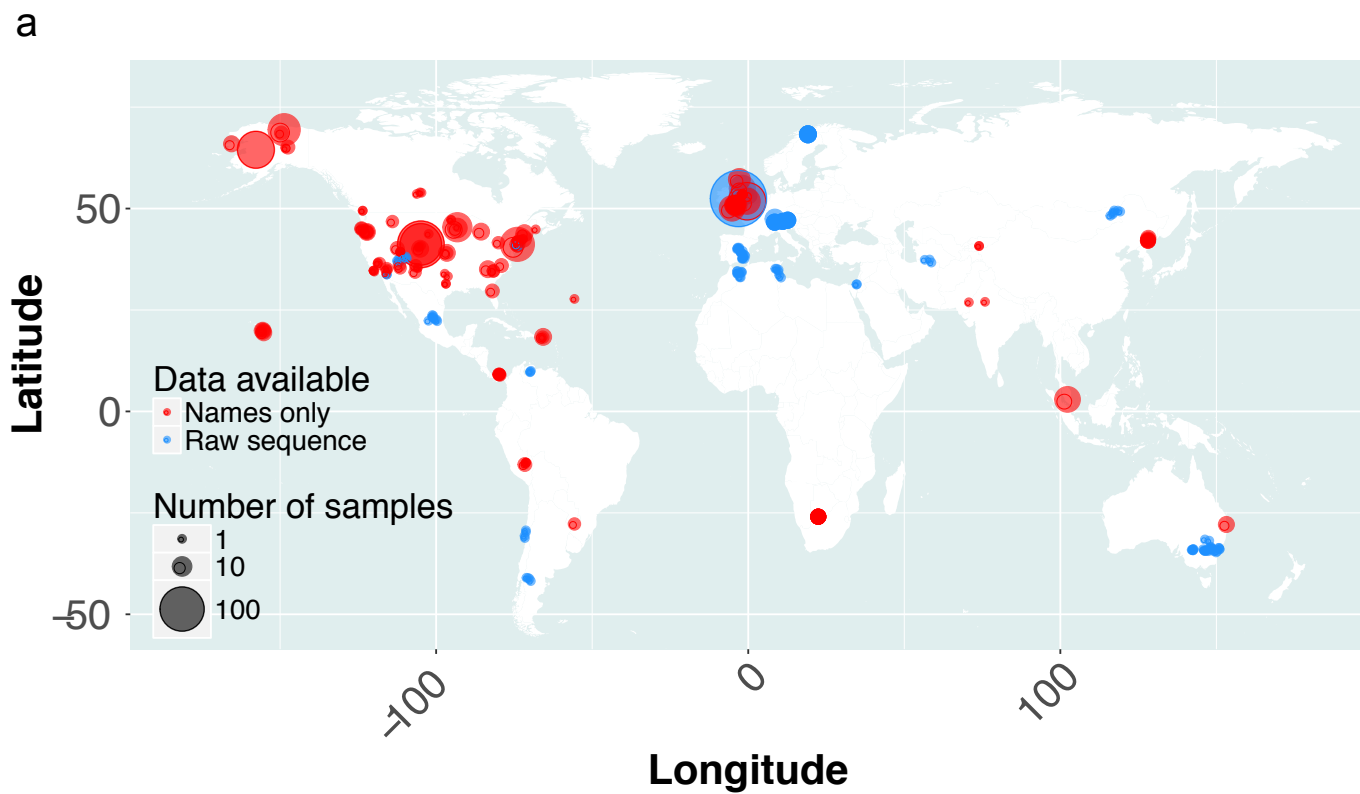
**Figure 4: The importance of bacterial taxa classified at different taxonomic ranks.**

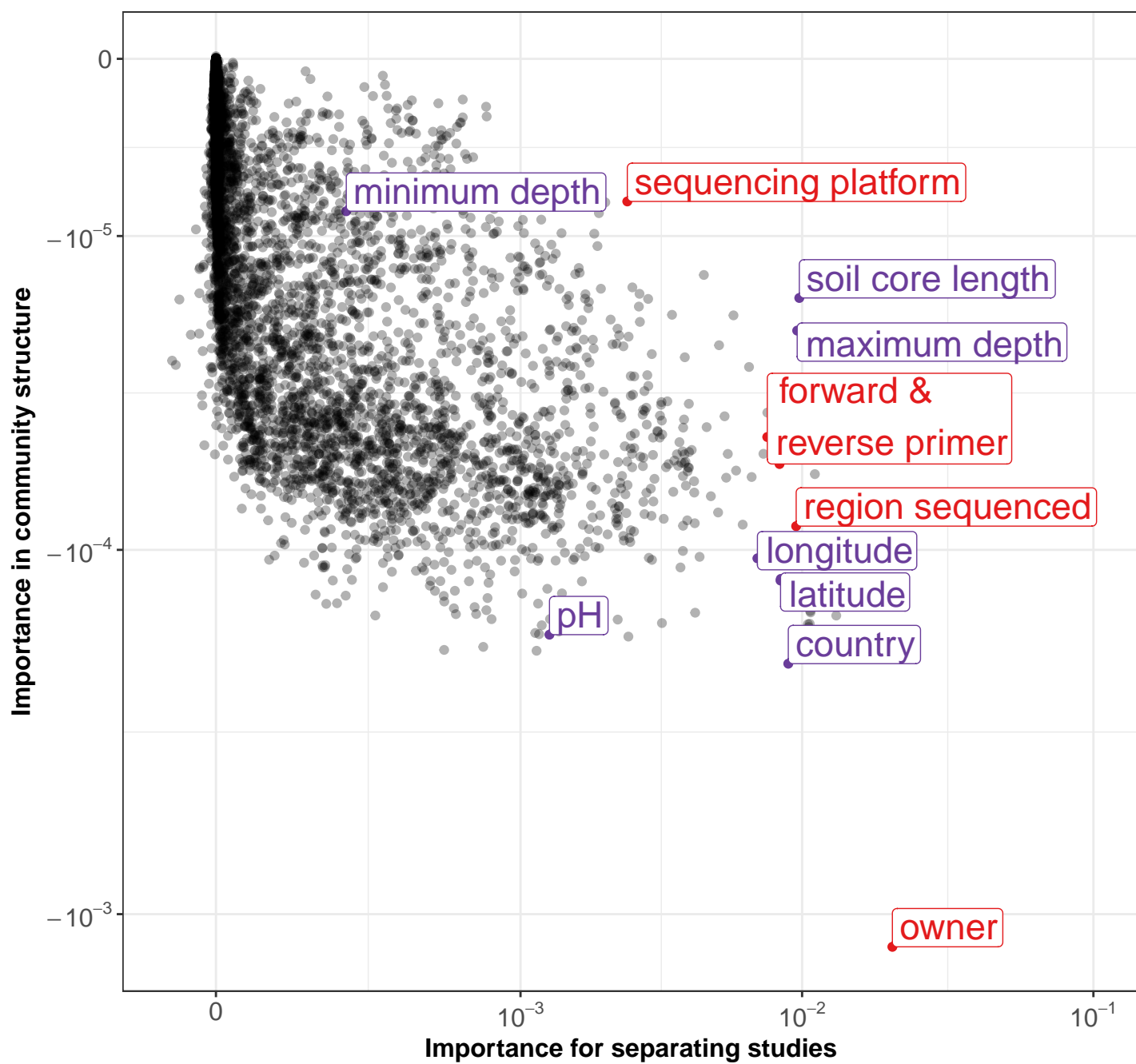
Lower taxonomic rank is more important for community structure (A), while high taxonomic rank is more important for separating studies (B). For each taxon, the difference was calculated between the variable importance (see *Methods*) of that taxon in a Random Forests model of either community structure or separating studies and the equivalent value from an analysis performed on the permuted dataset (see *Methods*). The lines and grey ribbons show the mean and standard error respectively of these values across taxa at each taxonomic rank considered.

**Figure 5: Importance of bacterial taxa in community structure related to their occurrence**

**in different studies.** The y-axis reports the variable importance in the Random Forests model of community structure (see *Methods*). Green ‘observed’ points correspond to those taxa shown in Figure 1. Orange ‘permuted’ points correspond to the same analysis on a null distribution (see *Methods*). Lines are general additive model (gam) smoothers. Each line is shown with a confidence interval (grey); where this is not visible it is narrower than the line it surrounds.



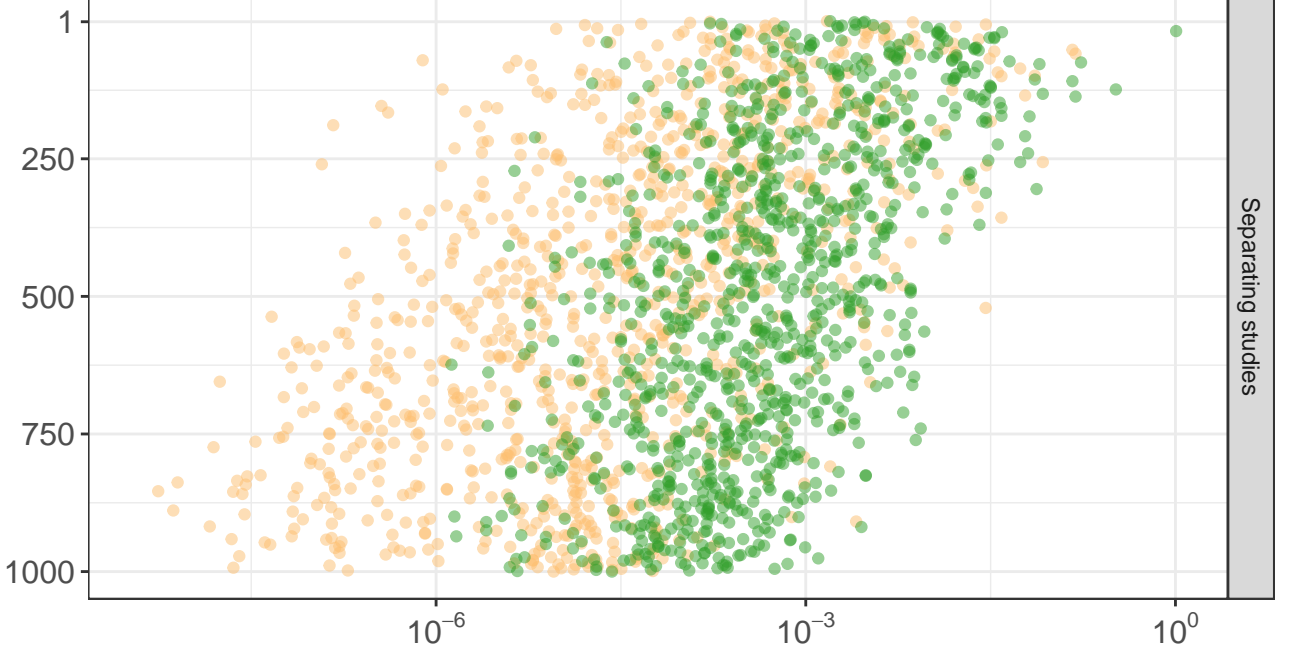




Rank importance

● observed  
● permuted

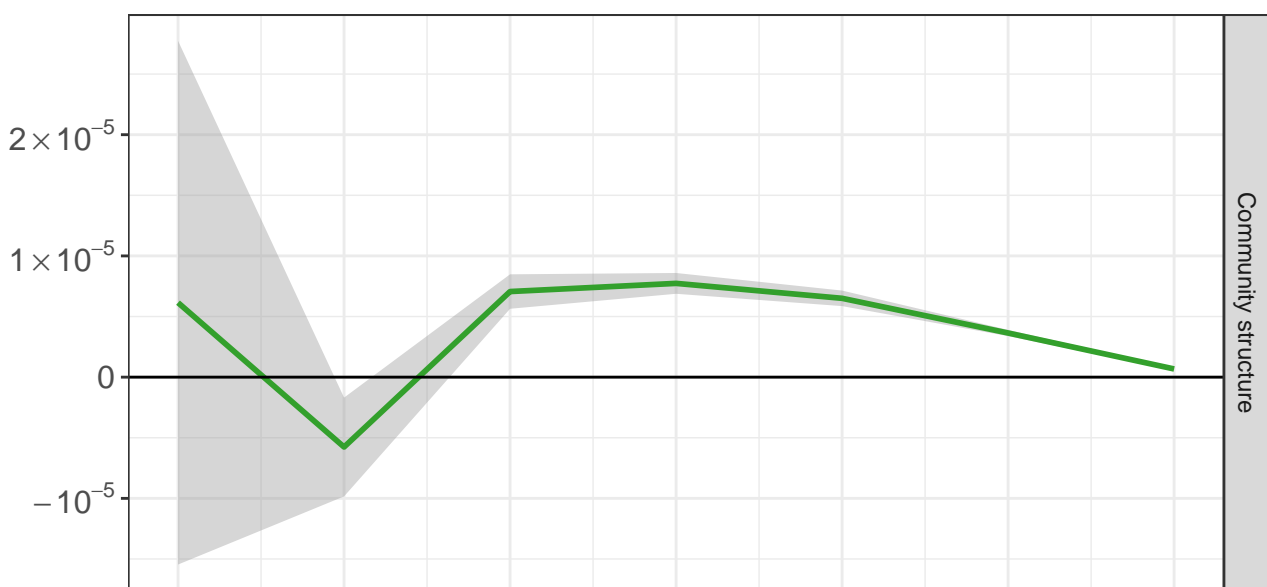
Community structure



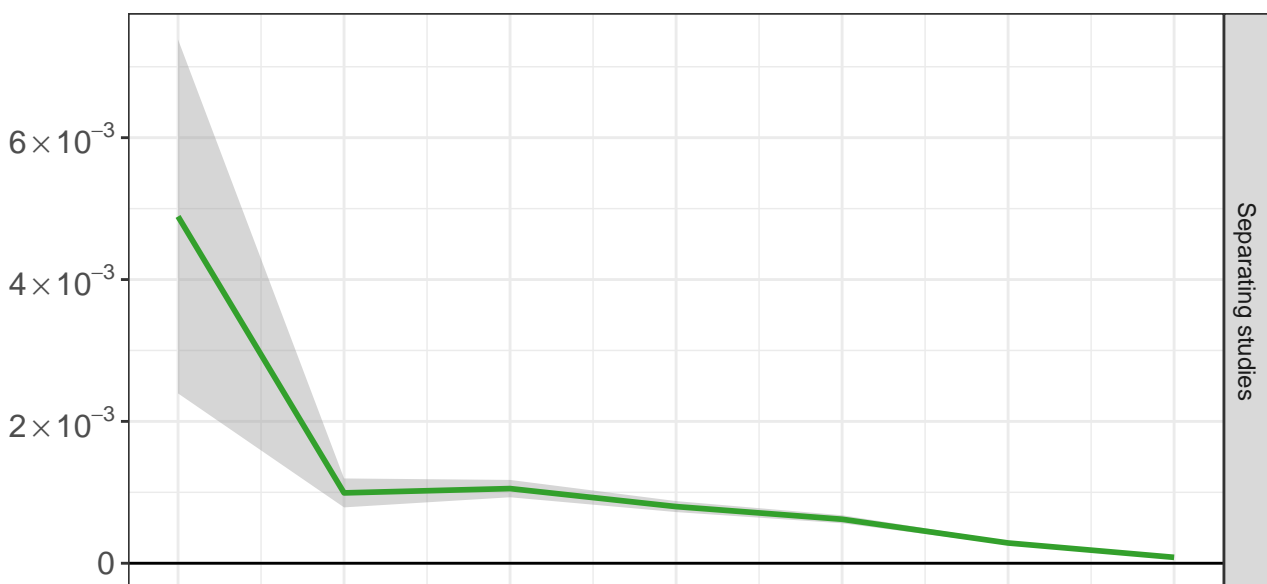
Separating studies

Average taxon abundance

Improvement on random



Community structure



Separating studies

Taxonomic level

Importance in community structure

Forest

— observed  
— permuted

Occurrence across studies %

0

$-2 \times 10^{-4}$

$-3 \times 10^{-4}$

0

50

100

species Kestanbolensis

genus Haemophilus

phylum Fusobacteria

genus Thauera

genus Brevundimonas

genus Thiobacillus

genus Ochrobactrum

family Streptococcaceae

family Brucellaceae

genus Sphingopyxis

genus Azospirillum

genus Methylothermobacter

order Lactobacillales

order Rhizobiales

phylum Actinobacteria

class Alphaproteobacteria

phylum Acidobacteria

class Actinobacteria

class Gemmatimonadetes

phylum Gemmatimonadetes

phylum Bacteroidetes